

09/26/00

09.27.00

A

JCS90 U.S. PAT. &amp; TM. OFF.

# UTILITY PATENT APPLICATION TRANSMITTAL

Attorney Docket No.	PD98-2384
First Inventor	Richard Fitzhugh Wrenn
Title	Method and Apparatus for Distributing Traffic Over Multiple Switched Fibre Channel Routes
Express Mail Label No.	EL700672297US

JCS90 U.S. PAT. &amp; TM. OFF.

09/26/00

## APPLICATION ELEMENTS

Assistant Commissioner for Patents  
Box Patent Application  
Washington, DC 20231

1. ☒ Fee Transmittal Form
2. ☒ Specification [ total pages 27 ]  
 - Descriptive title of the Invention  
 - Cross References to Related Applications  
 - Statement Regarding Fed sponsored R&D  
 - Reference to Microfiche Appendix  
 - Background of the Invention  
 - Brief Summary of the Invention  
 - Brief Description of the Drawings  
 - Detailed Description  
 - Claim(s)  
 - Abstract of the Disclosure
3. ☒ Drawings(g) [ total sheets 5 ]
4. ☒ Oath or Declaration [ unsigned or signed ]  
 [ total sheets 1 ]  
 a. ☒ Newly executed (original or copy)  
 b. ☐ Copy from prior appl. (37 C.F.R. § 1.63(d))  
 i. ☐ **DELETION OF INVENTOR(S)**  
 Signed statement attached deleting inventor(s) named in prior application, see 37 C.F.R. §§ 1.63(d)(2) and 1.33(b).

5. ☐ Microfiche Computer Program
6. Nucleotide/Amino Acid Sequence (if applicable)  
 a. ☐ Computer Readable Copy  
 b. ☐ Paper Copy (identical to computer copy)  
 c. ☐ Statement verifying identity of said copies

## ACCOMPANYING APPLICATION PARTS

7. ☒ Assignment Papers
8. ☐ 37 C.F.R. § 3.73(b) Statement (when there is an assignee) ☒ Power of Attorney
9. ☐ English Translation ☐ Copies of
10. ☐ IDS & Form 1449
11. ☐ Preliminary Amendment
12. ☒ Return Receipt Postcard (MPEP 503)
13. ☐ Small Entity ☐ Statement filed in prior application—Status still proper and desired
14. ☐ Certified Copy of Priority Document(s)
15. ☒ Other: Certificate of Mailing

\*NOTE FOR ITEMS 1-13: IN ORDER TO BE ENTITLED TO PAY SMALL ENTITY FEES, A SMALL ENTITY STATEMENT IS REQUIRED (37 C.F.R. § 1.22), EXCEPT IF ONE FILED IN A PRIOR APPLICATION IS RELIED UPON (37 C.F.R. § 1.29).

## 16. If a CONTINUING APPLICATION,

- ☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: 1  
 Prior application information: Examiner: \_\_\_\_\_ Group/Art Unit: \_\_\_\_\_

FOR CONTINUATION OR DIVISIONAL APPS ONLY: The entire disclosure of the prior application, from which an oath or declaration is supplied under Box 4b, is considered a part of the disclosure of the accompanying continuation or divisional application and is hereby incorporated by reference. The incorporation can only be relied upon when a portion has been inadvertently omitted from the submitted application parts.

## 17. CORRESPONDENCE ADDRESS

- ☒ Customer Number of Bar Code Label 25235 or ☒ Correspondence address below

Name	William J. KUBIDA, Esq.				
	Hogan & Hartson, LLP				
Address	1200 17 <sup>th</sup> Street Suite 1500				
City	Denver	State	CO	ZIP	80202
Country	US	Telephone	(719) 448-5900	Fax	(719) 448-5922

Name	Steven Kent Barton	Registration No.	36,445
(Signature)	<i>Steven K. Barton</i>	Date	25 Sept 2000

# FEE TRANSMITTAL for FY 2000

**TOTAL AMOUNT OF PAYMENT (\$)** **(\$730.00)**

## Complete if Known

Application Number	
Filing Date	herewith
First Named Inventor	Richard Fitzhugh Wrenn
Examiner Name	
Group / Art Unit	
Attorney Docket No.	PD98-2384

### METHOD OF PAYMENT (check one)

1. ☐ The Commissioner is hereby authorized to charge indicated fees and credit any over payments to:

Deposit Account Number **50-1123**

Deposit Account Name **Hogan & Hartson L.L.P.**

☒ Charge Any Additional Fee Required Under 37 CFR § 1.16 and 1.17

Payment Enclosed:

☒ Check ☐ Money Order ☐ Other

### FEE CALCULATION

#### BASIC FILING FEE

Entity Fee (\$)	Entity Fee (\$)	Fee Description	Fee Paid
690	345	Utility Filing Fee	<b>690.00</b>
310	155	Design filing fee	
480	240	Plant filing fee	
690	345	Reissue filing fee	
150	75	Provisional filing fee	

**SUBTOTAL (1) (\$690.00)**

#### 2. EXTRA CLAIM FEES

Total Claims	Extra Claims	Fee from below	Fee Paid
17	20**= 0	X	= 0
Independent Claims	3	-3**= 0	= 0
Multiple Dependent			= 0

\*\*or number previously paid, if greater; For Reissues, see below

#### Large Entity Small Entity

Fee Code (\$)	Fee Code (\$)	Fee Code (\$)	Fee Description
103 18	203 9		Claims in excess of 20
102 78	202 39		Independent claims in excess of 3
104 260	204 130		Multiple dependent claim, if not paid
109 78	209 39		**Reissue independent claims over original patent
110 18	210 9		**Reissue claims in excess of 20 and over original patent

**SUBTOTAL (2)**

**(\$)**

### FEE CALCULATION (continued)

Entity Fee (\$)	Entity Fee (\$)	Fee Description	Fee Paid
130	65	Surcharge - late filing fee or oath	
50	25	Surcharge - late provisional filing fee or cover sheet	
130	130	Non-English specification	
2,520	2,520	For filing a request for reexamination	
920*	920*	Requesting publication of SIR prior to Examiner action	
1,840*	1,840*	Requesting publication of SIR after Examiner action	
110	55	Extension for reply within first month	
380	190	Extension for reply within second month	
870	435	Extension for reply within third month	
1,360	680	Extension for reply within fourth month	
1,850	925	Extension for reply within fifth month	
300	150	Notice of Appeal	
300	150	Filing a brief in support of an appeal	
260	130	Request for oral hearing	
1,510	1,510	Petition to institute a public use proceeding	
110	55	Petition to revive - unavoidable	
1,210	605	Petition to revive - unintentional	
1,210	605	Utility issue fee (or reissue)	
430	215	Design issue fee	
580	290	Plant issue fee	
130	130	Petitions to the Commissioner	
50	50	Petitions related to provisional applications	
240	240	Submission of Information Disclosure Stmt	
40	40	Recording each patent assignment per property (times number of properties)	40 00
760	380	Filing a submission after final rejection (37 CFR § 1.129(a))	
760	380	For each additional invention to be examined (37 CFR § 1.129(b))	
Other fee (specify)			

\*Reduced by Basic Filing Fee Paid

**SUBTOTAL (3) (\$40.00)**

#### SUBMITTED BY

Name **Steven Kent Barton**

(Print/Type)

Signature *Steven K. Barton*

Registration No. (Attorney/Agent)

**36,445**

#### Complete (if applicable)

Telephone **(719) 448-5900**

Date **25 Sept 2000**

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of:

Richard Fitzhugh Wrenn

Serial No. \_\_\_\_\_

Filed: Herewith

For: Method and Apparatus for Distributing  
Traffic Over Multiple Switched Fibre  
Channel Routes

Group Art Unit: \_\_\_\_\_

Examiner: \_\_\_\_\_



CERTIFICATE OF MAILING BY EXPRESS MAIL

BOX PATENT APPLICATION

Assistant Commissioner of Patents and Trademarks  
Washington, D.C. 20231

Sir:

The undersigned hereby certifies that the following documents:

1. Utility Patent Application Transmittal;
2. Fee Transmittal and \$690 filing fee;
3. Utility Patent Application;
4. Executed Declaration;
5. Executed Power of Attorney by Assignee;
6. 5 sheets of drawings;
7. Recordation Form Cover Sheet PTO 1595 with Executed Assignment and Recording Fee of \$40.00;
8. Return postcard; and
9. Certificate of Mailing By Express Mail

relating to the above application, were deposited as "Express Mail", Mailing Label No. EL700672297US with the United States Postal Service, addressed to The Commissioner of Patents and Trademarks, Washington, D.C., 20231, Sept 26, 2000.

Sept. 26, 2000  
Date

Julie M. Trout  
Mailer

26 Sept 2000  
Date

Steven K. Barton  
Steven Kent Barton, Reg. No. 36,445  
Hogan & Hartson LLP  
One Tabor Center  
1200 17<sup>th</sup> Street, Suite 1500  
Denver, Colorado 80202  
Tel (719) 448-5906  
Fax (303) 899-7333

## **METHOD AND APPARATUS FOR DISTRIBUTING TRAFFIC OVER MULTIPLE SWITCHED FIBRE CHANNEL ROUTES**

### **FIELD OF THE INVENTION**

5 The invention relates to the field of computer networks. In particular, the invention relates to distributing network traffic between a pair of networked machines over multiple available routes through a network interconnecting the machines.

### **NATURE OF THE PROBLEM**

10 Most modern computer networks, including switched Fibre Channel networks, are packet oriented. In these networks, data transmitted between machines is divided into chunks of size no greater than a predetermined maximum. Each chunk is typically packaged with a header and a trailer into a packet for transmission. In Fibre Channel networks, packets  
15 are known as Frames.

20 Packets encounter delay while being routed through a network. Many networks have switches or routers that receive packets, store them, and forward the packets on towards their destinations when communications resources become available; storing and forwarding of packets introduces delay. Additional delay may be caused by propagation delay in the network interconnect between machines or switches of the network.

25 The multiple packets, or frames, associated with a single Fibre Channel operation are known as a

exchange. A Sequence is a group of one or more frames, forming part of an exchange, transmitted in a single direction over the network. A sequence may contain data, status, or control information. Each exchange may contain one or more sequences, and may contain data sequences of multiple frames with control and acknowledgment sequences that are often single frames. A Fibre Channel network having at least one switch is a switched Fibre Channel fabric. A Fibre Channel switch is a routing device generally capable of receiving frames, storing them, decoding destination information from headers, and forwarding them to their destination or to another switch further along a path toward their destination.

A network interface of a switch for connection of the switch to a machine is known as an F\_Port. An F\_Port having the ability to connect to a Fibre Channel Arbitrated Loop is known as an FL\_Port. An E\_Port is a network interface of a switch for connection of that switch to another switch of a fabric. A G\_Port is a port having the ability to operate as either an F\_Port or an E\_Port; and a GL\_Port further has the ability to connect to a Fibre Channel Arbitrated Loop. For purposes of this patent F\_Port includes any port of a switch that connects through a link to a machine, whether it be an F\_Port, G\_Port, GL\_Port, or an FL\_Port. Further, for purposes of this patent, an E\_Port includes any port of a switch that connects through a link to another switch, regardless of whether it be an E\_Port, GL\_Port, or G\_Port. Further, for purposes of this

patent, the term switch port includes any port of a switch, whether it be an E\_port or F\_port as defined herein.

5 A network interface for connection of a machine to a Fibre Channel fabric is known as an N\_Port, and a machine attached to a Fibre Channel network is known as a node. An L\_Port is a network interface for connection of a machine to a Fibre Channel Arbitrated Loop, and an NL\_Port is an N\_Port also  
10 having the ability to connect to a Fibre Channel Arbitrated Loop. For purposes of this patent, the term N\_Port includes both N\_Ports and NL\_Ports.

Machines, or "Nodes", attached to a Fibre Channel network may be computers, or may be storage devices  
15 such as RAID systems, disk drives, or other storage servers.

A Fibre Channel exchange operates between an originator N\_Port and a responder N\_Port. For example, an originator N\_Port may request an I/O  
20 operation such as a disk write; the machine attached to the responder N\_Port performs the operation. N\_Ports may be originators for some exchanges, and responders for others. Each Fibre Channel N\_Port is assigned identification for use as a destination  
25 address for frames intended for it, this identification is unique to the specific Fibre Channel network at a given time. Each Fibre Channel N\_Port participating in an exchange assigns exchange identification to that exchange, that exchange  
30 identification being unique among the exchanges in

progress on that N\_Port but not necessarily unique across the network.

For purposes of this application, a link is the data transmission and reception hardware and any associated firmware that form a connection between an N\_Port and an F\_Port of a switch, or between E\_Ports of two switches, of a Fibre Channel fabric. A link may incorporate a Fibre Channel Arbitrated Loop.

In a computer network, there may be more than one possible path, or sequence of links, switches, hubs, routers, etc. that may be traversed by a frame, between two machines attached to the network. Multiple paths may be intentional, providing extra capacity or redundant paths to protect against switch, node, or line failures, or may be unintentional consequences of network topology. Multiple paths between a pair of N\_Ports may exist if there are two or more switches in the network.

It is known that frames routed on different paths through a network may suffer different delays. Further, delay on each path varies with traffic on each link of the path, the arbitration sequence of each arbitrated loop forming part of a link, flow control delays like those often injected to avoid buffer overflow, and switch loading.

Machines transmitting data on modern high-speed networks usually do not wait for each frame to be acknowledged before transmitting following frames - multiple frames of a single Fibre Channel sequence may exist in a Fibre Channel fabric at the same time.

Further, frames of multiple sequences of a single exchange may also exist simultaneously in a Fibre Channel fabric, as may frames of multiple exchanges originated by any given N\_Port.

5           If frames of a sequence are transmitted on different paths through a fabric, an early-transmitted frame suffering long delay on one path may arrive at its destination after a late-transmitted frame that suffers little delay on  
10 another path. Frames transmitted on different paths thus may arrive at the destination N\_Port out-of-order, meaning that they are received in a different order than they were transmitted by their originating machine.

15           Frames received out-of-order may, and often do, require collection and sorting into correct order before they can be fully processed by the receiving machine. Some network protocols, including the TCP Internet protocol, presume out-of-order delivery and  
20 require that receiving machines collect and re-order frames before executing any command associated with them. Other order-dependent protocols, including the FCP protocol for encapsulating the SCSI storage interface protocol over Fibre Channel, assume that  
25 frames arrive in correct order - requiring that the Fibre Channel fabric deliver frames in-order. Some order-dependent protocols detect, and permit retry of, out-of-order frames even if they do not require that destinations perform resequencing. Fibre  
30 Channel frame headers include a sequence count field



with which out-of-order frames may be detected within a sequence.

Fibre Channel fabrics support a variety of order-dependent and order-independent protocols running on  
5 top of their low-level Fibre Channel mechanism.

Since frames transmitted over the same path through a network tend to arrive in order, many Fibre Channel systems permitting order-dependent protocols restrict communication between any two N\_Ports to  
10 transmission over one active path in each direction. Any other path between the N\_Ports may be usable as an alternate path should an active path fail, but may remain little used until that failure occurs. Networks that failover from an active path to an  
15 alternate path are known in the art of Fibre Channel networks. Frame routing of this type is known herein as static routing with alternate paths.

Links of an active path, especially links between switches, may be shared with traffic between other  
20 N\_Ports, including N\_Ports of other machines. As loads and network configurations change, it is possible for a statically routed active path to become a bottleneck while alternate paths may have unused capacity. It is desirable to make use of any  
25 available, otherwise unused, capacity of these alternate paths to provide improved network throughput.

It is known that many machines, including RAID storage subsystems, have the ability to queue  
30 multiple commands for execution. For example, a RAID

system may queue several read or write commands, received from one or more machines. Once queued, these commands are executed from the queue to or from cache, or to or from disk, in an order depending on availability of data in cache, disk availability and disk rotation. With proper interlocks, execution may often be in an order different from that in which the commands were received.

Commands that may be queued in these devices may include commands from multiple processes, or threads, running on a single machine having one or more processors. For example, a transaction-processing system may have several processes running, each process requiring access to a different record of a database on a RAID system, all requesting access to the database at about the same time. Each process may then create read, write, lock, or unlock commands for the database. Queuing and execution of each of these commands requires that an exchange of frames be transmitted between the machine and the device.

Fibre channel frame headers have a D\_ID field that encodes identification of the destination N\_Port of the frame. They also have an S\_ID field that encodes identification of the originating port of the frame. There is also an OX\_ID field that encodes the exchange identifier assigned by the originating N\_Port, and an RX\_ID field that encodes the exchange identifier assigned by the receiving N\_Port of the exchange. Since the receiving N\_Port does not assign RX\_ID until the exchange has begun and a frame is sent in response to other frames of the exchange, the

RX\_ID field of early frames of an exchange, including the first frame sent by the originating N\_Port, may not match the RX\_ID of late frames of the exchange.

### SOLUTION TO THE PROBLEM

5           A network, such as a Fibre Channel fabric, having two or more machines attached, each attached to the fabric through at least one N\_Port, has a first and a second path between an N\_Port of a first machine and an N\_Port of a second machine. The first machine  
10 originates several commands for execution on the second machine and embeds those commands and associated data in frames. Frames belonging to a first command are recognized and transmitted between the first and second machines over the first path,  
15 while frames belonging to a second command are transmitted between the first and second machines over the second path.

          Frames belonging to an individual exchange are recognized through the OX\_ID field of the frame  
20 headers. In an alternative embodiment, frames belonging to an individual exchange are recognized through a combination of the OX\_ID and the S\_ID fields of the frame headers. These fields, together with the destination address (D\_ID) of the frame, are  
25 input to a function whose output is used by routing and distributing tasks of one or more switches to index routing tables at a switch of the network fabric. These routing tables contain information determining the link over which each frame will be  
30 sent through the fabric from that switch towards the

destination. In this way, the routing tables determine paths, from what may be a multiplicity of possible paths, that each frame will follow through the network.

5 Except when routing tables are being updated, frames relating to the same exchange therefore follow the same path through the network, and therefore arrive in-order. Frames of simultaneous, but different, exchanges may be routed over different  
10 paths thus distributing traffic between the available paths.

As nodes, switches, and links are added to or removed from the network, and as a load-balancer adjusts demand on elements of the network, the  
15 routing tables are updated to reflect valid paths through the network and desired frame distribution among them. If more than one valid path appears in the table for any given destination, commands to that destination will tend to be distributed between the  
20 paths according to the frequency with which each path appears in the table.

## BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 is an illustration of a Fibre Channel network having several machines and several paths  
25 between some of these nodes;

Figure 1A, an illustration of multiple processes causing overlapping exchanges on an N\_Port;

Figure 2A, an example of frames for a simple write exchange;

Figure 2B, an example of frames for a simple read exchange;

Figure 3, an illustration of a Fibre Channel frame, as known in the art, detailing header information associated with the frame;

Figure 3A, an illustration of a prior-art routing table for routing frames based upon D\_ID;

Figure 3B, an illustration of a prior-art routing table for routing frames based upon S\_ID and D\_ID;

Figure 3C, an illustration of a routing table of the present invention for routing frames based upon D\_ID and OX\_ID;

Figure 3D, an illustration of a routing table of the present invention for routing frames based upon D\_ID, OX\_ID, and S\_ID;

Figure 4A, an illustration of a routing table system incorporating separate D\_ID and OX\_ID hash functions ahead of a routing table; and

Figure 4B, an illustration of a routing table system incorporating separate D\_ID and OX\_ID hash functions ahead of, and a level of indirection after, a base routing table;

## **DETAILED DESCRIPTION OF THE ILLUSTRATED EMBODIMENT**

A switched Fibre Channel network (Figure 1) has at least two machines, with a switched Fibre Channel fabric 100 interconnecting them. The fabric may incorporate two or more switches.

Machines of the network may include computers 102  
104, and 120, and RAID or other storage systems 106  
each having at least one N\_Port 108, 112, 114, 118,  
and 122, for interconnection to the fabric. Each  
5 N\_Port 108, 112, 114, 118, and 122 connects through a  
link 130, 134, 136, 140, and 142 to a switch of the  
switches 150, 152, and 154 of the fabric 100.  
Switches 150, 152, and 154 of the fabric may further  
be interconnected by additional links 160, 162, and  
10 164. Switches of the fabric may be joined by  
multiple links, switch 152 connects to switch 154 by  
a redundant link 165.

There may be, and preferably are, more than one  
path between a first and a second machine of the  
15 network. There are frequently also more than one  
possible path from a first N\_Port to a second N\_Port.  
For example, computer 120 may communicate to RAID  
system 106 through a first path comprising N\_Port  
122, link 142, switch 150, link 162, switch 154, link  
20 140, and N\_Port 118; or through a second path  
comprising N\_Port 122, link 142, switch 150, link  
160, switch 152, link 164, switch 154, link 140, and  
N\_Port 118. A third path may also exist similar to  
the second path but using the redundant link 165 from  
25 switch 152 to switch 154, comprising N\_Port 122, link  
142, switch 150, link 160, switch 152, link 165,  
switch 154, link 140, and N\_Port 118. Similarly,  
computer 102 may communicate with computer 104  
through a path comprising N\_Port 108, link 130,  
30 switch 150, link 162, switch 154, link 136 and N\_Port  
114, or through an alternative path comprising N\_Port

108, link 130, switch 150, link 160, switch 152, link 164, switch 154, link 136, and N\_Port 114.

Consider the first and second path described above between computer 120 and RAID system 106. In a network utilizing static routing, only one of these paths is active at a given time. The active path may include one or more elements that become overloaded, or become a bottleneck for these communications. For example, if the active path from N\_Port 108 of computer 102 to N\_Port 114 of computer 104 is through link 162 and the active path from N\_Port 122 of computer 120 to N\_Port 118 of RAID system 106 is also through link 162, it is possible for link 162 to have a heavy load while link 160 is idle.

There may be multiple processes simultaneously executing on computer 120. Each of these processes 200 and 202 (Figure 1A) may generate an I/O request 204 and 206 as known in the art, each of which in turn is performed through an exchange 208 and 210 as known in the art. These exchanges may overlap in time as they are transferred by the N\_Port 122 to and from the fabric; overlapping I/O operations may result from multiple concurrent processes on a machine and many other known causes. For example but not by way of limitation, a disk write operation and a disk read operation may overlap.

A disk-write command may be packetized as a write exchange Figure 2A comprising a write command frame 250 sent from the originating N\_Port 251 to a receiving N\_Port 252, and a write-data sequence 254

sent after a transfer ready frame 255 is received by  
the originating N\_Port 251. When writing to cache or  
disk has been completed by the receiving N\_Port's  
machine, a response status frame 256 is returned to  
5 the originating N\_Port 251. Additional  
acknowledgment and control frames may also be used.  
Similarly, a disk-read I/O command becomes a read  
exchange, Figure 2B, which operates through  
transmission of at least a read command frame 260  
10 from the originating N\_Port 251 to a receiving N\_Port  
252, which may be associated with a RAID system or  
other storage device. When data associated with the  
read operation is ready, the receiving N\_Port 252  
returns a data sequence 264 and status 266 frames to  
15 the originating N\_Port 251, which may be associated  
with a computer. The write exchange of Figure 2A may  
overlap the read exchange of Figure 2B. For example  
and not by way of limitation, it is possible that the  
originating port read command 260 may be transmitted  
20 by the originating port 251 after the write command  
frame 250 is transmitted and before the transfer  
ready frame 255 is received by the originating port  
251.

Each frame, or packet, transmitted over a Fibre  
25 Channel network has structure as illustrated in  
Figure 3. The frame contains a header, an optional  
payload, and a trailer. The header includes several  
fields, including a Destination Identification (D\_ID)  
field 300, a Source Identification (S\_ID) field 302,  
30 an Originator Exchange Identifier (OX\_ID) 304, and a  
Responder Exchange Identifier (RX\_ID) 306. The RX\_ID



306 may change during an exchange because it is assigned by the responder node after the first frames of an exchange are received by that node; the OX\_ID 304 is stable within an exchange. It is possible for a switch to nearly-simultaneously receive frames having identical D\_ID 300 and OX\_ID 304 fields from different sources, having different S\_ID fields 302.

A switch of a Switched Fibre Channel Fabric receives frames having the format of Figure 3, and typically has multiple switch ports, such as E\_Ports 170 and 178 (Figure 1), and F\_Ports 174 and 176 of switch 150. Once the switch 150 receives a frame on an incoming switch port it is expected to forward that frame on a selected outgoing port of the switch. The selected outgoing port is a switch port, other than the incoming switch port, on a path from the originating N\_Port to the receiving N\_Port.

It is known that a routing table 330, Figure 3A indexed by a hash function 332 of the D\_ID 300 field of a frame header, may be used to generate an outgoing port selector for controlling the outgoing switch port on which frames are forwarded by the switch. The D\_ID 300 field is transformed by a hash-function 332 to an address 334, the address locating a table entry in the table 330. Each entry has an outgoing port selector 336 that controls the switch port on which the frame is forwarded by the switch.

In an effort to improve the ability of network management software to optimize traffic flow on a network, some switches input the S\_ID field 302

5 (Figure 3B) of the frame, or an incoming switch port number on which the frame was received, to a hash function 342 in addition to the D\_ID field 300. As in the routing system of Figure 3A, the hash function 342 generates an address 344 that locates a table entry in a routing table 346. The table entry then provides an outgoing port selector 348. This permits the switch to route traffic to a given destination from two different sources over two different routes.

10 In a switch of the present invention, a routing table 350, Figure 3C, is indexed by an address 354 generated by a hash function 352 of the D\_ID field 300 and the OX\_ID field 304 of each frame header. An outgoing port selector 356 is derived from a table entry of the routing table 350 located in the table by the address 354. The outgoing port selector 356 is used to control the switch port on which frames are transmitted.

20 In an alternative embodiment of a switch of the present invention, the S\_ID field 302, as well as the D\_ID field 300 and the OX\_ID field 304, of each frame header is used by a hash function 360 (Figure 3D) to generate an address 362. Address 362 is then used to generate an outgoing port selector 364 by reading a table entry from a routing table 366. This embodiment provides opportunity to independently control frame distribution between available paths for each source.

30 Consider frames received by a switch 150 of the present invention from computer 120 and intended for

RAID system 106 N\_Port 118. The headers of each of these frames are decoded by switch 150. In the network as illustrated, frames having D\_ID field 300 corresponding to a destination of N\_Port 118 may reach that destination through a path through switches 152 and 154, and through a second path through switch 154 directly. A hash function of the D\_ID field 300 and at least one bit of the OX\_ID field 304 of the header are therefore used to index routing table 180 to select the outgoing switch port. The routing table 180 has the structure illustrated in Figures 4C or 4D. The hash function is selected such that all entries of the routing table 180 that may be selected by a valid D\_ID field 300 correspond to a valid outgoing port on a path to the N\_Port identified by D\_ID that is distinct from the incoming switch port.

Frames belonging to the same exchange have the same OX\_ID field; therefore these frames follow the same route through the network and tend to arrive in-order within that exchange. Frames may, however, arrive out-of-order with respect to frames of other exchanges.

In a Fibre Channel network, there may be paths between two ports that are "better" in some way than others. Multiple bits of the OX\_ID field 304 may be considered by a routing table to distribute frames between a preferred and a less preferred path. For example, if three bits of OX\_ID are considered by a routing table of switch 150, eight table entries may be addressed for the same D\_ID. If three of these

have an outgoing port selector specifying E\_Port 170,  
while five specify E\_Port 178, about three-eighths of  
frames will tend to follow the path through switches  
150 and 154 while five-eighths of frames will tend to  
follow the path through switches 150, 152, and 154.  
If more than one valid path appears in the table for  
any given destination, exchanges directed to that  
destination are thus distributed between the paths  
according to the frequency with which each path  
appears in the table.

As machines, switches, and links are added to or  
removed from the network the routing tables are  
updated to reflect valid paths through the network  
and the desired frame distribution among them. The  
routing tables are also adjusted as a load-balancer  
task, which may run on any compute-capable machine or  
switch of the network, adjusts demand on elements of  
the network. For example, should the link 162  
attached to E\_Port 170 of switch 150 fail, those  
routing table entries specifying this port may be  
replaced by entries specifying E\_Port 178 so that  
frames may reach their intended destination.

It is not necessary that the hash function 340  
consider all bits of the OX\_ID field, it is expected  
that significant distribution of traffic among  
multiple routes can be achieved by considering as few  
as one or several low bits of the OX\_ID field.

In an alternative embodiment of the present  
invention, a hash function 400 (Figure 4A) of the  
D\_ID field 300 generates an address-X 402 for a two-

dimensional routing table 404. A second hash function 406 generates an address-Y 408 for the routing table 404 from the OX\_ID field 304 and may also consider the S\_ID field 302. The routing table generates a outgoing port selector 410 as previously described. The routing table 404 therefore has a predetermined, number of port entries for each valid D\_ID, each entry of which is readily locatable. The set of port entries for a particular D\_ID are referenced as a line of the routing table.

The embodiment of Figure 4A is advantageous because only one line of the routing table need be rewritten to alter the distribution of frames between paths to an individual N\_Port. Further, this embodiment lends itself to control of frame distribution among paths because the number of entries associated with each destination is constant and these entries are readily located in the table.

While the routing table of the present invention has been described as producing an outgoing port selector from a hash function of the D\_ID and OX\_ID fields 300 and 304, that operation may be either direct or indirect. In an alternative embodiment, a level of indirection is used such that paths may be taken in or out of service quickly, without need to rewrite many of the outgoing port selectors in the routing table. For example, consider the routing table structure of Figure 4B. In this embodiment, a hash function 420 of the D\_ID field 300 generates an address-X 422. A second hash function 424 of at least one bit of the OX\_ID field 304, and,

optionally, the S\_ID field 302, produces an address-Y 426. The address-X 422 and the address-Y are combined to address a routing table 428. The routing table 428 thereupon produces a path code 430. Path code 430 is then translated by a portmap table 432 into the outgoing port selector 434. Path code 430 may have more bits than outgoing port selector 434.

In this embodiment, should a link fail it may be possible to rewrite the portmap table 432 to reroute all frames onto a functioning link (if one exists) in less time than it would take to restructure the routing table 428. Once the frames are rerouted onto a functioning link by rewriting the portmap table 432, the routing table 428 may be adjusted to balance the load. Alternatively, if path code 430 has more bits than the outgoing port selector 434, it may not be necessary to rewrite the routing table 428.

Routing tables of the present invention may be implemented in firmware or hardware of the switch. It is known that implementation of routing tables in hardware provides advantage for switches having heavy load and large numbers of switch ports. In a hardware implementation, routing table 350 of Figure 3C, 366 of Figure 3D, 404 of Figure 4A, or 428 of Figure 4B, may be implemented with a static RAM, and the portmap table 432 with a second static RAM. In such an embodiment, the routing table address inputs are multiplexed so it can be written by a processor of the switch such that the processor can maintain the routing table. The routing table is thereby addressable by either the address generated by the

hash function or functions, or by an address generated by the processor.

5 The hash function used for addressing the routing table may be any of many hash functions known in the art of computer science. This function may also comprise concatenation of a preselected group of bits of each input to the hash function; such as concatenation of one or more low-order OX\_ID bits with several bits of the D\_ID field to produce an index to the routing table. This function may also comprise concatenation of functions of bits from each field, or concatenation of bits of results of a function applied to each field.

10 A computer program product is a machine-readable memory having recorded on it a program for performing a particular function; this may be a read-only memory or may be an erasable and rewritable memory such as RAM, CD-RW, tape, flash memory, or magnetic disk. It is anticipated that routing control software for controlling routing tables as herein described may be distributed or operated as a computer program product. Similarly, a switch containing firmware for constructing and utilizing the routing table of the present invention in routing frames is expected to contain memory having that firmware, and therefore contains a computer program product.

20 While much reference has been made to a first and second path through the network, the invention is not limited to a pair of paths. The invention is

applicable to any reasonable number of concurrently available paths between nodes of a network.

While the invention has been particularly shown and described with reference to a preferred embodiment thereof, it will be understood by those skilled in the art that various other changes in the form and details may be made without departing from the spirit and scope of the invention.



What is claimed is:

1. A network comprising:

a network fabric further comprising at least two switches and a plurality of links, each link connected to at least one switch of the at least one switch;

a first N\_Port connected to a link of the network fabric;

a second N\_Port connected to a link of the network fabric;

wherein there exists a first path and a second path from the first N\_Port to the second N\_Port through the network fabric;

wherein network traffic from the first N\_Port to the second N\_Port is automatically distributed between the first path and the second path by the switch such that frames transmitted in a first direction and related to any single exchange are transmitted over the same path of the first and second path yet frames transmitted in the first direction and related to different but overlapping exchanges need not follow the same path.

2. The network of Claim 1, wherein the frames related to the any single exchange are identified by a switch as belonging to the single exchange through fields of a frame header comprising an originator exchange identifier field.

3. The network of Claim 2, wherein frames are routed by at least one routing table located within a switch of the at least two switches, the routing table having inputs comprising a hash function of a

5 destination identifier of the frame header and at  
6 least one bit of the originator exchange identifier.

1 4. The network of Claim 3, wherein the network  
2 comprises a switched Fibre Channel fabric.

1 5. The network of Claim 4, wherein the hash function  
2 has input further comprising a field selected from  
3 the group consisting of a source identifier field of  
4 the frame header and an incoming port number on which  
5 the frame was received by the switch.

1 6. The network of Claim 4, wherein the routing table  
2 produces an index to a second table that provides an  
3 outgoing port identifier for the switch.

1 7. The network of Claim 4 wherein a load-balancing  
2 task of the network updates the at least one routing  
3 table to alter a distribution of exchanges among  
4 paths.

1 8. The network of Claim 4, wherein the hash function  
2 further comprises a first hash sub-function of at  
3 least one bit of the destination identifier and  
4 having an output, a second hash sub-function of the  
5 at least one bit of the originator exchange  
6 identifier and having an output, and a concatenation  
7 operation of the output of the first hash sub-  
8 function with the output of the second hash sub-  
9 function.

1 9. The network of Claim 8, wherein the second hash  
2 sub-function is a bit select operation.

1 10. The network of Claim 8, wherein the hash function  
2 has inputs further comprising an input selected from  
3 the group consisting of an incoming port identifier  
4 on which the frame was received and at least one bit  
5 of a source identifier field of the frame header.

1 11. A program product for distributing network  
2 traffic between a first N\_Port of a network and a  
3 second N\_Port of a network, the network having a  
4 plurality of paths for frames from the first N\_Port  
5 to the second N\_Port and at least one switch, the  
6 program product operable upon said switch and  
7 comprising computer-readable code for:

8 maintaining a routing table, the routing table  
9 indexed by an output of a hash function of inputs  
10 comprising a destination identification field and an  
11 originator exchange identifier field of a header of a  
12 frame;

13 causing the routing table to be accessed upon  
14 receipt of a frame, the routing table coupled to  
15 determine a selected port for transmission of the  
16 frame; and

17 causing the frame to be transmitted on the  
18 selected port.

1 12. The program product of Claim 11, wherein the hash  
2 function has inputs further comprising a an input  
3 selected from the group consisting of a source  
4 identifier field of the frame header and an identity  
5 of a switch port upon which the frame was received.

1 13. The program product of Claim 11, wherein the  
2 routing table is coupled to determine a selected port

3 by providing an index to a second table that provides  
4 a selected port identifier.

1 14. A switch for a network capable of distributing  
2 frames received on a first port over a plurality of  
3 ports, the switch comprising

4 a plurality of ports including a first port, the  
5 first port capable of receiving a frame;

6 a routing table capable of determining a port of  
7 the plurality of ports for forwarding a received  
8 frame based upon an address;

9 a hash function generator capable of generating  
10 an address for the routing table based upon  
11 information comprising a destination identification  
12 field and at least one bit of an originator exchange  
13 identifier field of a header of the received frame;

14 a processor for maintaining the routing table;  
15 and

16 apparatus for receiving a frame and for passing a  
17 received frame to the port determined by the routing  
18 table.

1 15. The switch of Claim 14 wherein the hash function  
2 generator is capable of generating an address for the  
3 routing table based upon information further  
4 comprising an identifier selected from the group  
5 consisting of a source identifier field of the header  
6 of the received frame and an port identifier of the  
7 switch port on which the frame is received.

1 16. The switch of Claim 14, wherein the hash function  
2 generator further comprises devices to perform the  
3 hash function of a destination identification field

4 and at least one bit of an originator exchange  
5 identifier field of the header of the received frame,  
6 and the routing table comprises a memory capable of  
7 being addressed by the address generated by the hash  
8 function.

1 17. The switch of Claim 16, wherein the memory of the  
2 routing table is implemented by at least one RAM, the  
3 RAM being writable by the processor and coupled to be  
4 addressed through a multiplexor capable of selecting  
5 a RAM address from the group of addresses comprising  
6 an address generated by the processor and the address  
7 generated by the hash function.

## ABSTRACT

A computer network has two or more switches and a plurality of links. A first machine and a second machine are interconnected by the network in such a way that there exist multiple paths through the network from an N\_Port of the first machine to an N\_Port of the second machine. Network traffic from the N\_Port of the first machine to the N\_Port of the second machine is distributed between the multiple paths such that frames related to any single exchange traverse the same path yet frames of a first exchange need not traverse the same path as frames of a second exchange. Frames of each exchange therefore tend to be received by their destination in order with respect to other frames of that exchange, while they are not necessarily received in-order with respect to frames of other exchanges.

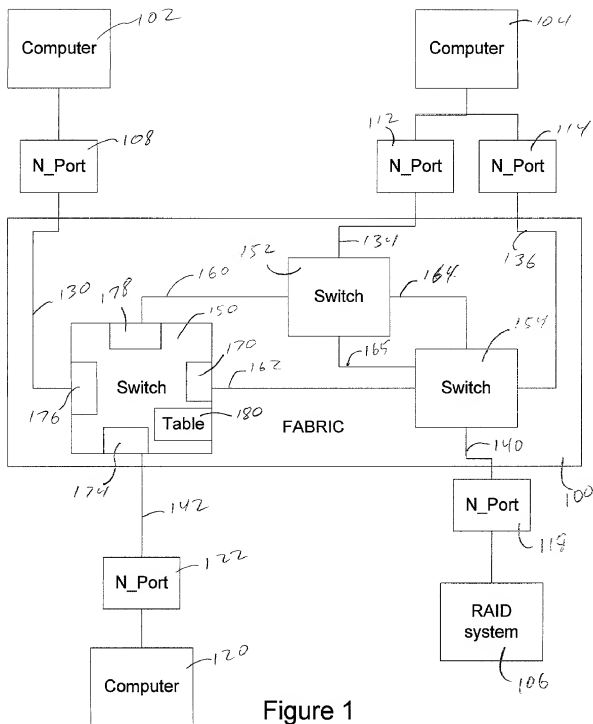
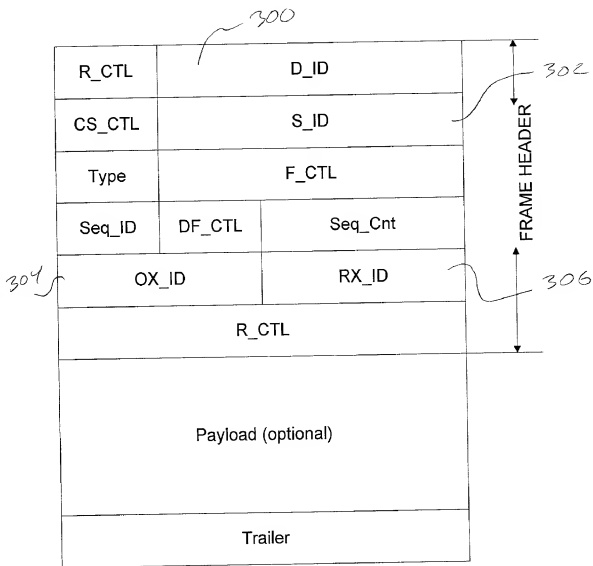


Figure 1





2025 RELEASE UNDER E.O. 14176



Packet Structure

Figure 3

Figure 1 is a block diagram labeled "Prior Art". It shows a data flow process. On the left, two input lines are labeled "D\_ID" and "S\_ID". These inputs feed into a block labeled "Hash Function". The output of the "Hash Function" is a line labeled "Address" with the handwritten value "342" next to it. This "Address" line feeds into a block labeled "Routing Table". The output of the "Routing Table" is a line that leads to an "Outgoing Port". The "Routing Table" block has a handwritten value "46" next to it, indicating the output size or type.

Figure 3D is a block diagram illustrating a routing process. It features three input fields on the left: D\_ID (labeled 340), S\_ID (labeled 342), and OX\_ID. These inputs feed into a central block labeled "Hash Function" (labeled 344). The output of the Hash Function is an "Address" (labeled 346), which is then fed into a "Routing Table" block (labeled 348). The final output of the Routing Table is the "Outgoing Port Selector" (labeled 350).

Figure 3D

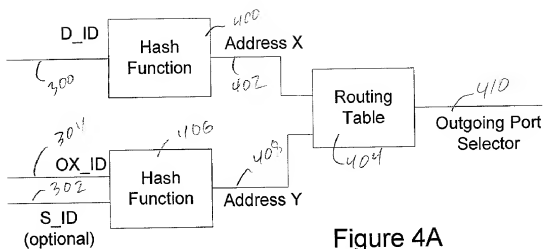


Figure 4A

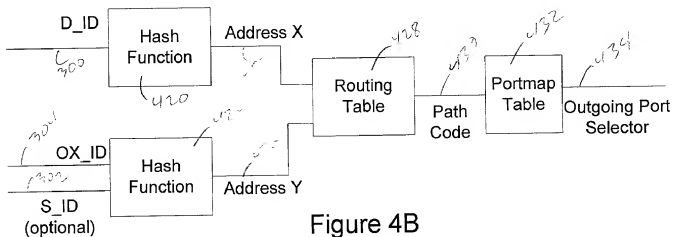


Figure 4B

DECLARATION

As a below named inventor, I hereby declare that: my residence, post office address, and citizenship are as stated below next to my name. I believe I am the original, first, and sole inventor (if only one name is listed below) or a joint inventor (if plural inventors are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled:

Method and Apparatus for Distributing Traffic Over Multiple Switched Fibre Channel Routes

as described in the specification ☒ attached or ☐ of patent Application Serial No.  
filed and amended on \_\_\_\_\_.

I hereby state that I have reviewed and understand the contents of the above-identified specification, including the claims, as amended by any amendment referred to above; that I do not know and do not believe the same was ever known or used in the United States of America before my or our invention thereof, or patented or described in any printed publication in any country before my or our invention thereof or more than one year prior to this application; that the invention has not been patented or made the subject of an inventor's certificate issued before the date of this application in any country foreign to the United States of America on an application filed by me or my legal representative or assigns more than twelve months prior to this application; and that I acknowledge the duty to disclose information of which I am aware which is material to the examination of this application in accordance with Title 37, Code of Federal Regulations  $\square$  1.56(a). Such information is material when it is not cumulative to information already of record or being made of record in the application, and


- ☐ (1) it establishes, by itself or in combination with other information, a prima facie case of unpatentability of a claim; or  
☐ (2) it refutes, or is inconsistent with, a position the applicant has taken or may take in:  
    (i) opposing an argument of unpatentability relied on by the Office, or  
    (ii) asserting an argument of unpatentability.

I hereby claim foreign priority benefits under Title 35, United States Code  $\square$  119 of any foreign application(s) for patent or inventor's certificates listed below and have also identified below any foreign application(s) having a filing date before that of the applications(s) on which priority is claimed:

COUNTRY	APPLICATION NUMBER	Date Filed	Priority Claimed under 35 USC 119
			<input type="checkbox"/> YES <input type="checkbox"/> NO

I hereby claim the benefit under Title 35 United States Code  $\square$  120 of any United States application(s) listed below and, insofar as any subject matter of any claim of this application is not disclosed in the prior United States Application, I acknowledge the duty to disclose material information as defined in Title 37, Code of Federal Regulations  $\square$  1.56(a) which occurred between the filing date of the prior application and the national PCT international filing date of this application.

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

FULL NAME OF SOLE OR FIRST INVENTOR		INVENTOR'S SIGNATURE	DATE
Richard Fitzhugh Wrenn			9-22-2000
RESIDENCE		CITIZENSHIP	
2360 Oak Hills Drive		USA	
Colorado Springs, CO 80919			
POST OFFICE ADDRESS			
Same			
FULL NAME OF SECOND JOINT INVENTOR		INVENTOR'S SIGNATURE	DATE
		CITIZENSHIP	
POST OFFICE ADDRESS			

SEP-25-00 11:43 From: HOGAN &amp; HARTSON

7194485922

T-011 P.02/02 Job-205

## IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

## Applicant/Patentee:

Richard Fitzhugh Wrenn

## Serial No.:

Date Filed: Herewith

For: Method and Apparatus for  
Distributing Traffic Over Multiple Switched  
Fibre Channel Routes

Attorney File No.: 68854.0123

Digital Docket No.: PD98-2384

**POWER OF ATTORNEY BY ASSIGNEE**

Under the provisions of 37 C.F.R. § 3.71, the undersigned assignee of record of the entire interest in the above-identified patent/patent application by virtue of an assignment recorded (check as applicable):



Concurrently Herewith



Date Recorded \_\_\_\_\_



Reel \_\_\_\_\_ Frame \_\_\_\_\_

elects to conduct the prosecution of the application/maintenance of the patent to the exclusion of the inventor(s). The undersigned hereby declares that she has reviewed the above-referenced assignment and hereby declares that, to the best of her knowledge, title is in the Assignee, and further declares that all statements made herein of her own knowledge are true and that all statements made on information and belief are believed to be true. The assignee hereby revokes any previous powers of attorney and appoints the following to prosecute this application/maintain this patent and transact all business in the Patent and Trademark Office connected therewith:

**(Prosecuting Attorney List)**

Irene Kosturakis, Reg. No. 33,724

Rich Lange, Reg. No. 27,295

Louis Brucculeri, Reg. No. 38,834

Sarah T. Harris, Reg. No. 35,891

Joseph Arrambide, Reg. No. 38,589

Keith Lutsch, Reg. No. 31,861

Theodore S. Park, Reg. No. 26,971

William J. Kubida, Reg. No. 29,864

Stuart T. Langley, Reg. No. 33,940

Carol W. Burton, Reg. No. 35,465

Steven Kent Barton, Reg. No. 36,445

E. Michael Byorick, Reg. No. 34,131

Matthew G. Dyor, Reg. No. 42,278

Steven C. Petersen, Reg. No. 36,238

Sarah S. O'Rourke, Reg. No. 41,226

Please direct all communications relative to this application to the following addressee:

WILLIAM J. KUBIDA  
Hogan & Hartson LLP  
One Tabor Center  
1200 17th Street, Suite 1500  
Denver, Colorado 80202  
(719) 448-5909

**ASSIGNEE****COMPAQ COMPUTER CORPORATION**Date: 25 September 2000BY: mboutsham

NAME: Diane H. Strang

TITLE: Administrator, Patents

Authorized To Sign This Document On Behalf Of Compaq Computer Corp.